

Eq. 1

Optimization

Partial derivatives:

Defn. $f: \mathbb{R}^d \rightarrow \mathbb{R}$, { e_i } standard basis, if

$$\lim_{h \rightarrow 0} \frac{f(a + he_i) - f(a)}{h} \text{ exists.}$$

we say f has partial derivative w.r.t. x_i at a

$$\frac{\partial f(a)}{\partial x_i} \text{ or } \frac{\partial}{\partial x_i} f(a)$$

Treat everything else as a constant and take the 1-d derivative

let $a \in \mathbb{R}^d$

$$\text{Eq. } f(x) = a^T x. \quad \frac{\partial f(x)}{\partial x_i} = \frac{\partial}{\partial x_i} \left(\sum_{j=1}^d a_j x_j \right) = \sum_{j=1}^d \frac{\partial}{\partial x_i} a_j x_j = a_i.$$

Gradient Fn. $f: \mathbb{R}^d \rightarrow \mathbb{R}$. Assume that partial derivatives of f all

Under some conditions exist in a neighborhood of x and are also continuous
 that will hold across at x . Then f is differentiable and its grad satisfies:
 this course.

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_d} \end{bmatrix}.$$

Convention

$$\frac{\partial f}{\partial x} = \nabla f(x)^T$$

Eg. let $a \in \mathbb{R}^d$ and $f(x) = a^T x$. $\nabla f(x) = \left[\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_d} \right]^T$

$$= [a_1, a_2, \dots, a_d]^T = a$$

Eg. let $W \in \mathbb{R}^{d \times d}$ and $f(x) = x^T W x$

$$\begin{bmatrix} x_1, \dots, x_d \end{bmatrix} \begin{bmatrix} w_{11} & w_{12}, \dots, w_{1d} \\ w_{21}, w_{22}, \dots, w_{2d} \\ \vdots & \ddots \\ w_{d1}, \dots, w_{dd} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}$$

$$\left[\sum_{j=1}^d x_j w_{ij}, \sum_{j=1}^d x_j w_{2j}, \dots, \sum_{j=1}^d x_j w_{dj} \right] \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$= \sum_{i=1}^d \sum_{j=1}^d x_i w_{ij} x_j = \sum_{i=1}^d w_{ii}^2 x_i + \sum_{i=1}^d \sum_{j \neq i} x_i w_{ij} x_j$$

$$\frac{\partial f(x)}{\partial x_m} = \sum_{i=1}^d \sum_{j \neq i} x_i w_{ij} x_j$$

$$= \sum_{i=1}^d \sum_{j \neq m} x_m w_{mj} x_j + \sum_{i \neq m} x_i \sum_{j \neq i} w_{ij} x_j$$

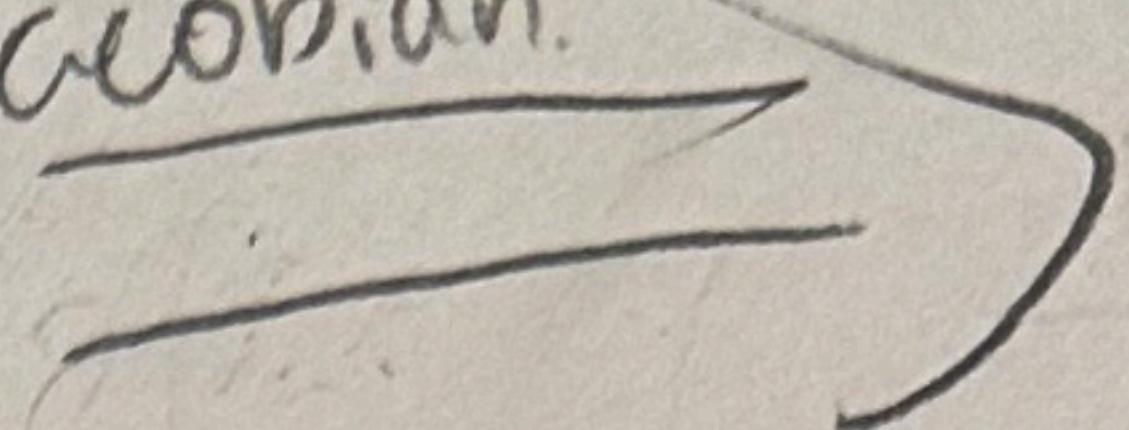
$$= \sum_{i=1}^d w_{mi} x_i + \sum_{i \neq m} x_i w_{im}$$

$$= \sum_{j=1}^d w_{mj} x_j + \sum_{i=1}^d x_i w_{im}$$

$$= (W[m, :])^T x + (W[:, m])^T x$$

$$= (W[m, :])^T x + (W^T[m, :])^T x$$

$$\begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_d} \end{bmatrix} = \begin{bmatrix} (w[1,:] + w^T[1,:])^T x \\ (w[2,:] + w^T[2,:])^T x \\ \vdots \\ (w[d,:] + w^T[d,:])^T x \end{bmatrix} = (w + w^T)x$$

Jacobian 

Optimization Algos

1 Line Search

for each step t

find a direction $u^{(t)}$

find $\alpha \in \mathbb{R}_+$ to minimize $h(\alpha) = f(w^{(t)} + \alpha u^{(t)})$

$$w^{(t+1)} = w^{(t)} - \alpha u^{(t)}$$

$u^{(t)}$ can be ^{arbitrary} coordinates \Rightarrow coordinate descent

If $u^{(t)}$ is the gradient direction \Rightarrow gradient descent

Newton's Method

$$w^{(t+1)} = w^{(t)} - \nabla f(w^{(t)})$$

Jacobian. A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$. $J_f(x) = \frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f_1}{\partial x} \\ \vdots \\ \frac{\partial f_m}{\partial x} \end{pmatrix} = \begin{pmatrix} \frac{\partial f_i}{\partial x_j} \end{pmatrix}_{ij}$

$\frac{\partial f_i}{\partial x} = \frac{\partial f(x)_i}{\partial x}$ are differentiable

Eg. $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ $f(x) = Ax$, $A \in \mathbb{R}^{k \times d}$. $Ax \in \mathbb{R}^k$

$$\frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial(Ax)_1}{\partial x} \\ \frac{\partial(Ax)_2}{\partial x} \\ \vdots \\ \frac{\partial(Ax)_k}{\partial x} \end{pmatrix} = \begin{pmatrix} \frac{\partial(A[1,:])x}{\partial x} \\ \frac{\partial(A[2,:])x}{\partial x} \\ \vdots \\ \frac{\partial(A[k,:])x}{\partial x} \end{pmatrix} = \begin{pmatrix} A[1,:] \\ A[2,:] \\ \vdots \\ A[k,:] \end{pmatrix} = A$$

Chain rule. $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$.

$g \circ f: \mathbb{R}^n \rightarrow \mathbb{R}^p$

f, g differentiable.

$$\text{Proof: } [J_{g \circ f}(x)]_{ik} = \frac{\partial(g \circ f)_i}{\partial x_k} = \sum_{j=1}^m \frac{\partial g_i(f(x))}{\partial f_j(x)} \cdot \frac{\partial f_j(x)}{\partial x_k}$$

$\frac{\partial g_i(f(x))}{\partial x_k}$

$$= \frac{\partial g(f(x)_1, f(x)_2, \dots, f(x)_m)}{\partial x_k}$$

$\frac{\partial f_j(x)}{\partial x_k}$

$\frac{\partial g(f(x))}{\partial f(x)}$

$\frac{\partial f(x)}{\partial x_k}$

$$J_{g \circ f}(x) = J_g(f(x)) J_f(x)$$

$$= \sum_{j=1}^m (J_g(f(x)))_{ij} (J_f(x))_{jk}$$

$$= [J_g(f(x)) J_f(x)]_{ik}$$

Hessian twice differentiable

$f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{pmatrix}$$

Usually, $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$
If 2nd order deriv. are then symmetric.

$$f(x) = x^T A x$$

$$\nabla f(x) = (A + A^T)x = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{pmatrix}$$

$$\frac{\partial(\frac{\partial f}{\partial x_i})}{\partial x_j} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

$$= A[j,:] + A^T[j,:]$$

Taylor's theorem

$$\nabla^2 f(x) = A + A^T$$

Taylor's theorem (Multivariate case)

$f: D \rightarrow \mathbb{R}$ where $D \subseteq \mathbb{R}^n$ is a convex set

$$f(x) = f(u) + \nabla f(u)^T (x-u) + (x-u)^T \nabla^2 f(\bar{s})(x-u) + o(\|x-u\|^2)$$

Go back to Optimization Algs.

Newton's method.

$$\bar{f}(w) = f(w^{(t)}) + \nabla f(w^{(t)})^T (w - w^{(t)}) + \frac{1}{2} (w - w^{(t)})^T \nabla^2 f(w^{(t)}) (w - w^{(t)})$$

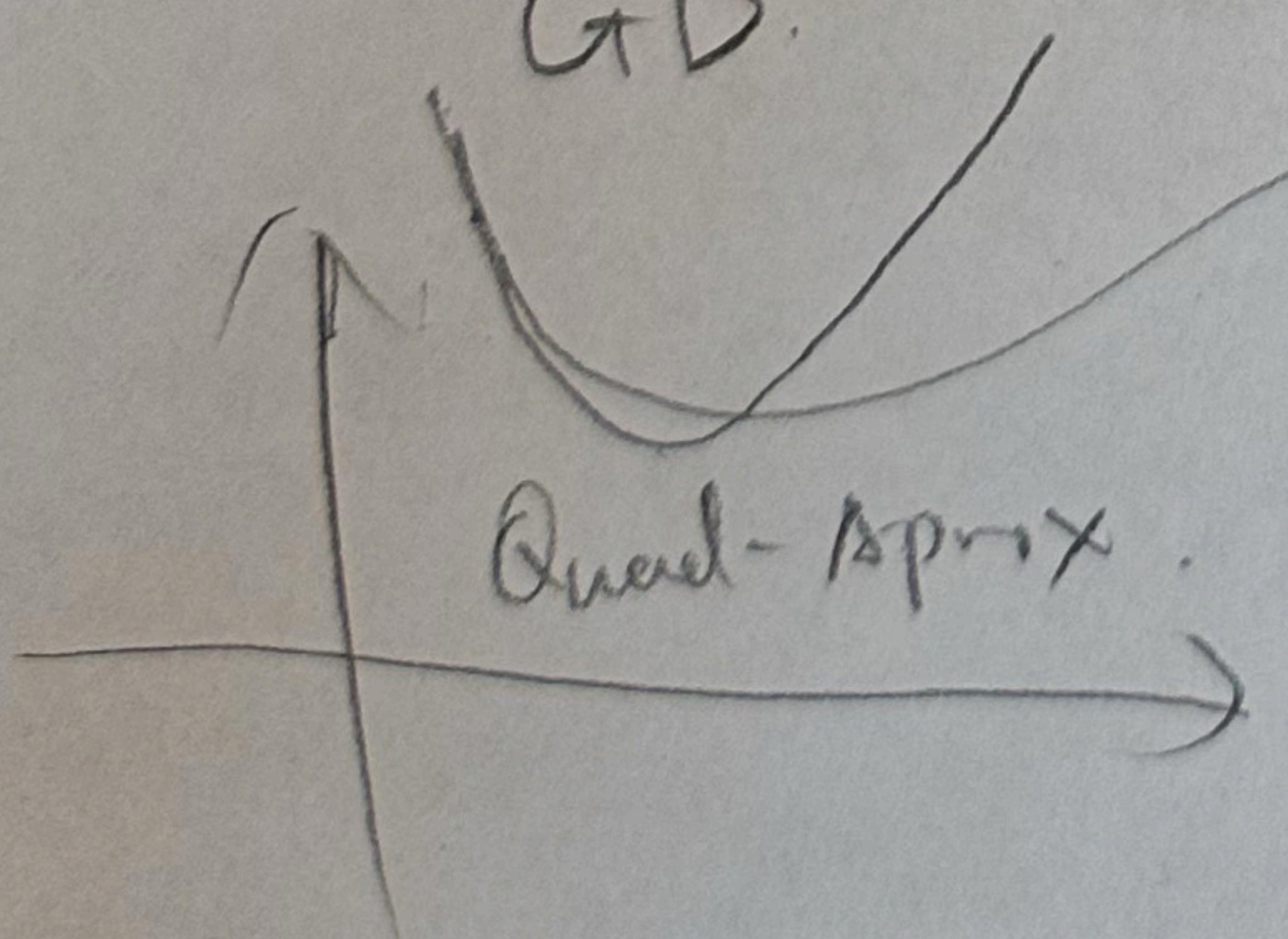
Quadratic function

then minimize $\bar{f}(w) \Rightarrow$ take derivative, set to 0.

$$\nabla \bar{f}(w) = \nabla f(w^{(t)})^T + 2 \cdot \frac{1}{2} \nabla^2 f(w^{(t)}) (w - w^{(t)}) \underset{\text{set}}{=} 0$$

$$\text{Newton } w^{(t+1)} = w^{(t)} - (\nabla^2 f(w^{(t)}))^{-1} \nabla f(w^{(t)})$$

$$\text{GD. } w^{(t+1)} = w^{(t)} - \alpha_t \nabla f(w^{(t)})$$



Another view of Newton's method 2-d case.

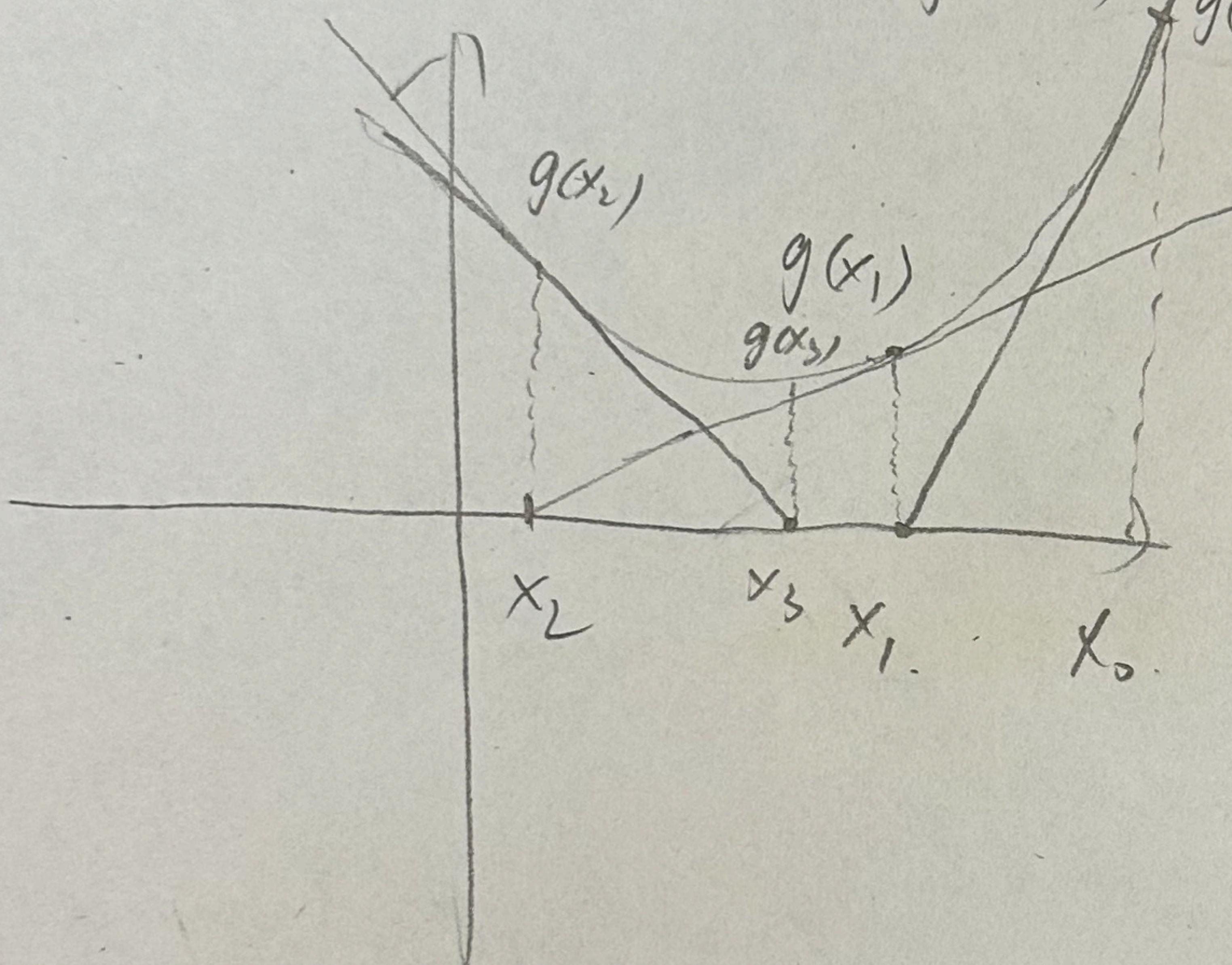
Want to find the critical pt., that is, find \hat{w} s.t.

$$f(\hat{w}) = 0$$

let $g(w) = f'(w)$ linear approximation

$g(w) \approx g(w_0) + f'(w_0)(w - w_0) \Rightarrow$ find the zero-point of
this lin-approx.

$$w_1 = w_0 - \frac{f(w_0)}{f'(w_0)} = w_0 - \frac{g'(w_0)}{g''(w_0)}$$



$$\text{Eq } l(\theta) = \sum_{i=1}^n y_i \log(\mu(x_i^\top \theta)) + (1-y_i) \log(1-\mu(x_i^\top \theta))$$

$$\nabla l(\theta) = \sum_{i=1}^n y_i \frac{1}{\mu(x_i^\top \theta)} \cdot \mu(x_i^\top \theta)^{x_i} + (1-y_i) \frac{1}{1-\mu(x_i^\top \theta)} \cdot (-\mu(x_i^\top \theta)) \cdot x_i$$

$$= \frac{1}{n} \sum_{i=1}^n y_i (1-\mu(x_i^\top \theta)) x_i + (1-y_i) (-\mu(x_i^\top \theta)) x_i$$

$$= -\frac{1}{n} \sum_{i=1}^n (y_i - \mu(x_i^\top \theta)) x_i$$

$$\nabla^2 \ln(\theta) = \frac{1}{n} \sum_{i=1}^n \mu(x_i^\top \theta) x_i x_i^\top$$

$$\mu(x_i^\top \theta) > 0. \quad \nabla^2 \ln(\theta) \succeq 0.$$

$\frac{1}{n} \mu(I - M)$

$$\theta^{(t+1)} = \theta^{(t)} - \left(\frac{1}{n} \sum_{i=1}^n \mu(x_i^\top \theta) x_i x_i^\top \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n (\mu(x_i^\top \theta) - y_i) \cdot x_i \right)$$

Pros of Newton's method.

For certain fns. ~~NR~~ converges ^{way} faster than GD

$$NR. \quad |f(x^{(t)}) - f(x_*)| \leq \underbrace{\left(\frac{1}{2}\right)^{2^{-t}}}_{t \in \mathcal{O}(\log(\log \frac{1}{\varepsilon}))} \downarrow$$

$$GD. \quad |f(x^{(t)}) - f(x_*)| \leq O(\frac{1}{t}), \quad t \approx O(\frac{1}{\varepsilon}).$$

Extreme example. If you have a Quadratic fn.

Newton's method converges in 1 step.

Cons of N-R.

3. sensitive to initial points. \rightarrow need good initialization

4. calculating inverse takes $O(d^3)$ operations \Rightarrow single step slow.

2. unpredictable: even w/ very nice functions, e.g.
 & Lipschitz
 it can still diverge.

Damped newton. $\text{W}_{t+1} = \text{W}_t - \gamma_t (\nabla^2 f(\text{w}_t))^{-1} \nabla f(\text{w}_t)$

1 Hessian not invertible \Rightarrow regularization.